

A fast algorithm for detecting maximal number of matched pairs under a given caliper

Pavel S. Ruzankin^{*1}

¹Sobolev Institute of Mathematics, Novosibirsk, Russia

¹Novosibirsk State University, Novosibirsk, Russia

Abstract

We present a new algorithm which detects the maximal number of matched disjoint pairs satisfying a given caliper when the matching is done with respect to a scalar index (e.g., propensity score), and constructs a corresponding matching. If each of the groups is ordered with respect to the index then the number of operations needed is $O(N)$, where N is the total number of objects to be matched. The case of 1-to- n matching is also considered.

Keywords: propensity score matching, matching with caliper.

1 One to one matching

We consider matching disjoint pairs of objects from two groups, which we will call, using common terminology, treated and control objects. In other words, a control object can be matched to no more than one treated object and vice versa. We will consider only one-dimensional distance, such as in propensity score matching, when the distance between objects is the distance between points on the real line corresponding to these objects, assuming each object is somehow projected to a unique point on the real line. We will call these points propensity scores of the objects for the sake of clarity. However no assumptions are made on how these points are related to the objects.

Let X_j , $j = 1, \dots, K$, and Y_j , $j = 1, \dots, L$, be the propensity scores of treated and control objects, K and L being the total numbers of treated and control objects, respectively. X_j and Y_j may take any values on the real line,

^{*}email: ruzankin@math.nsc.ru

not necessarily from $(0, 1)$. Let $N = K + L$. Let $c > 0$ be the caliper for our matching, i.e., we match only such i - j pairs that $|X_i - Y_j| \leq c$.

A natural problem for this setting is to find the maximal number of pairs that can be matched. Though this problem can be solved employing network flow optimization algorithms, the known algorithms have complexity not less than $O(N^2)$ (if no assumptions on sparsity of the matching are made). This approach to matching problems was used, e.g., by Rosenbaum (2012) and Pimentel et al. (2015).

Our main goal is to introduce a fast algorithm for detecting the maximal number of matched pairs and for constructing a corresponding matching. The presented algorithm has complexity $O(N)$ when both the treated and the control objects are sorted with respect to propensity score:

$$X_1 \leq X_2 \leq \dots \leq X_K \quad \text{and} \quad Y_1 \leq Y_2 \leq \dots \leq Y_L. \quad (1)$$

Thus once we have sorted the observations (which takes $O(N \log N)$ operations), we can reasonably fast solve the inverse problem of finding the minimal caliper suitable for using Q percent of data for a given Q . For instance, if propensity score belongs to the interval $(0, 1)$ then k iterations of the algorithm ($O(kN)$ operations) yield the accuracy of 2^{-k} for the minimal caliper.

From now on we assume that relation (1) holds.

The variable M will contain the current number of matched pairs. After the algorithm finishes, M contains the maximal number of matched pairs. A_k and B_k store the index numbers of control and treated object, respectively, in the k -th matched pair.

We present the algorithm as the following pseudocode:

```

 $M := 0$ 
 $i := 1$ 
 $j := 1$ 
while ( $i \leq K$  and  $j \leq L$ )
  if ( $|X_i - Y_j| \leq c$ )
     $M := M + 1$ 
     $A_M := i$ 
     $B_M := j$ 
     $i := i + 1$ 
     $j := j + 1$ 
  else
    if ( $X_i < Y_j$ )
       $i := i + 1$ 
    else
       $j := j + 1$ 
    end if
  end if
end while

```

As we see, the algorithm just walks through all the observations and successively collects all feasible pairs.

The algorithm requires $O(N)$ operations since in each iteration of the while-loop the variable i or j or both are increased. Certainly, to apply the algorithm, first we must sort the observations with respect to propensity score, which requires $O(N \log N)$ operations.

In Sec. 3 we prove that the algorithm produces the maximal possible number of matched pairs.

2 1-to- n matching

The algorithm can be modified for 1-to- n matching. We assume that a treated object is to be matched with no more than n control objects, and a control object must not be matched to more than one treated object. Our algorithm maximizes the number of matched control objects or, in other words, the number of matched pairs.

The following pseudocode uses the same variables as above. D_j is the number of controls matched to the j -th treated object.

```

 $M := 0$ 
 $i := 1$ 
 $j := 1$ 
 $k := 1$ 
 $D_j := 0$  for all  $j = 1, \dots, L$ 
while ( $i \leq K$  and  $j \leq L$ )
  if ( $|X_i - Y_j| \leq c$ )
     $M := M + 1$ 
     $A_M := i$ 
     $B_M := j$ 
     $D_j := k$ 
     $i := i + 1$ 
    if ( $k < n$ )
       $k := k + 1$ 
    else
       $k := 1$ 
       $j := j + 1$ 
    end if
  else
    if ( $X_i < Y_j$ )
       $i := i + 1$ 
    else
       $k := 1$ 
       $j := j + 1$ 
    end if
  end if
end while

```

The complexity is still $O(N)$ and does not depend on n since, as above, in each iteration of the while-loop the variable i or j or both are increased.

3 Validity of the algorithm

We offer the following two proofs for the algorithms above.

3.1 The first proof

First consider *one to one matching*. We will prove that the former algorithm yields the maximal number of matched pairs by induction.

There exists a matching \mathcal{M} satisfying the caliper c (i.e., $|X_i - Y_j| \leq c$ for all $(i, j) \in \mathcal{M}$) and containing the maximal number of matched pairs.

First steps of the algorithm skip the observations that can not be used for matching, i.e. treated objects i such that $X_i < \min_j Y_j - c$ and controls j such that $Y_j < \min_i X_i - c$.

After the above operation we can assume that $|X_1 - Y_1| \leq c$. Let us show that matching now the first treated with the first control object, as the algorithm does, does not reduce the maximal number of matched pairs, if we match the maximal number of pairs for the remaining $2, \dots, K$ -th treated and $2, \dots, L$ -th control objects.

If the first treated or the first control object are not matched in \mathcal{M} then removing from \mathcal{M} a possible pair with the first treated or the first control object and then adding $(1, 1)$ to \mathcal{M} does not change the number of pairs in \mathcal{M} . Thus, in this case, matching the pair $(1, 1)$ and then matching the maximal number of pairs for the $2, \dots, K$ -th treated and $2, \dots, L$ -th control objects yields the total maximal number of matched pairs.

The case when \mathcal{M} contains the pair $(1, 1)$ is clear.

It remains to consider the case when \mathcal{M} contains some pairs $(1, j_1)$ and $(i_1, 1)$, where $i_1 \neq 1$ and $j_1 \neq 1$. In this case we have $X_1 \leq X_{i_1} \leq Y_1 + c$ and $Y_1 \leq Y_{j_1} \leq X_1 + c$. Therefore $X_{i_1} - Y_{j_1} \leq Y_1 + c - Y_1 = c$ and $Y_{j_1} - X_{i_1} \leq X_1 + c - X_1 = c$, and, hence,

$$|X_{i_1} - Y_{j_1}| \leq c.$$

Thus removing from \mathcal{M} the pairs $(1, j_1)$ and $(i_1, 1)$ and adding the pairs $(1, 1)$ and (i_1, j_1) does not change the number of pairs in \mathcal{M} . Again, matching the pair $(1, 1)$ and then matching the maximal number of pairs for the $2, \dots, K$ -th treated and $2, \dots, L$ -th control objects yields the total maximal number of matched pairs.

Applying the above argument to the remaining $2, \dots, K$ -th treated and $2, \dots, L$ -th control observations proves the validity of the former algorithm by induction.

For the case of *one to n matching* it suffices to consider the second algorithm as the first one applied to observations where we take n identical treated objects instead of each corresponding treated object from the original observations, i.e., we “repeat” each treated object n times.

3.2 The second proof

For the second proof we use a theorem on a Monge-Kantorovich mass transfer problem.

Let P and Q be finite (nonnegative) continuous measures on the real line with Borel σ -algebra with $P(\mathbb{R}) = Q(\mathbb{R})$. Let

$$\rho(P, Q) = \inf_U \{U(\{(x, y) : |x - y| > c\})\},$$

where the supremum is taken over all (nonnegative) continuous measures on \mathbb{R}^2 with marginals P and Q : $U(A, \mathbb{R}) = P(A)$, $U(\mathbb{R}, A) = Q(A)$ for all Borel A . Put

$$F(t) = \mathbf{P}((-\infty, t)), \quad G(t) = \mathbf{P}((-\infty, t)).$$

Ruzankin (2001) proved the following theorem.

Theorem 1 *The equalities hold:*

$$\begin{aligned} \rho(P, Q) &= \lim_{y \rightarrow \infty} S(y) - P(\mathbb{R}) \\ &= \lim_{y \rightarrow \infty} T(y) - Q(\mathbb{R}), \end{aligned}$$

where the functions S and T are specified by the relations

$$\lim_{y \rightarrow -\infty} S(y) = \lim_{y \rightarrow -\infty} T(y) = 0, \quad (2)$$

$$dS(y) = \max\{dF(y), T(y + dy - c) - S(y)\}, \quad (3)$$

$$dT(y) = \max\{dG(y), S(y + dy - c) - T(y)\} \quad (4)$$

for all y and the assumption that the functions S and T are left continuous. The functions S and T exist and are uniquely defined.

The relations (3), (4) and the left-continuity condition mean that

$$\begin{aligned} S(y + w) &= S(y) + P([y, y + w)) \\ &\quad + \sup_{0 < v \leq w} (T(y + v - z) - S(y) - P([y, y + v)))^+, \end{aligned} \quad (5)$$

$$\begin{aligned} T(y + w) &= T(y) + Q([y, y + w)) \\ &\quad + \sup_{0 < v \leq w} (S(y + v - z) - T(y) - Q([y, y + v)))^+ \end{aligned} \quad (6)$$

for all y and $w > 0$, where $t^+ = \max\{t, 0\}$.

Put

$$\mu(P, Q) = \sup_U \{U(\{(x, y) : |x - y| \leq c\})\}, \quad (7)$$

where the supremum is taken over all continuous measures on \mathbb{R}^2 with marginals P and Q . We have

$$\mu(P, Q) = P(\mathbb{R}) - \rho(P, Q) = 2P(\mathbb{R}) - \lim_{y \rightarrow \infty} S(y) = 2Q(\mathbb{R}) - \lim_{y \rightarrow \infty} T(y).$$

The measure U_0 that achieves the supremum in (7) can be built as follows (Ruzankin 2001). Put

$$V(y) = F(y) - T(y + c) + G(y + c), \quad (8)$$

$$W(y) = G(y) - S(y + c) + F(y + c). \quad (9)$$

The functions $V, W, F - V, G - W$ are left continuous and nondecreasing.

Let the measure Z on \mathbb{R}^2 be defined by

$$Z((-\infty, x) \times (-\infty, y)) = \min\{V(x), W(y)\}. \quad (10)$$

Let the measure R be an arbitrary measure with marginals $F - V$ and $G - W$ (here we use a function of y instead of the corresponding measure of $(-\infty, y)$). Put $U_0 = Z + R$. Then the measure U_0 achieves the supremum in (7).

Here the measure Z is responsible for an optimal “mass transfer”:

$$Z(\{(x, y) : |x - y| \leq c\}) = Z(\mathbb{R}^2) = \mu(P, Q)$$

since $V(y - c) \leq W(y) \leq V(y + c)$ for all y .

To prove the optimality of the above algorithms we are to consider the case of discrete P and Q with supports consisting of finite numbers of points.

First consider one to one matching. Let the measure \tilde{P} be concentrated on the points $X_1 \leq X_2 \leq \dots \leq X_K$ and the measure \tilde{Q} be concentrated on the points $Y_1 \leq Y_2 \leq \dots \leq Y_L$ with

$$\tilde{P}(\{y\}) = \#\{j : X_j = y\}, \quad \tilde{Q}(\{y\}) = \#\{j : Y_j = y\},$$

where the $\#$ sign denotes the number of elements of a set.

If $K \neq L$ then to use Theorem 1 we have to extend one of the measures \tilde{P} or \tilde{Q} . Put

$$\begin{aligned} D &= \max\{X_K, Y_L\} + 2c, \\ P(A) &= \tilde{P}(A) + (Q(\mathbb{R}) - P(\mathbb{R}))^+ I_D(A), \\ Q(A) &= \tilde{Q}(A) + (P(\mathbb{R}) - Q(\mathbb{R}))^+ I_D(A), \end{aligned}$$

where $I_D(A) = 1$ if $D \in A$ and is zero otherwise. Now $P(\mathbb{R}) = Q(\mathbb{R})$ and we can apply Theorem 1.

The function S can increase only at the points $X_j, Y_j + c, D + 2c, D + 3c$ while the function T can increase only at the points $Y_j, X_j + c, D + 2c, D + 3c$. Relations (5), (6) can be rewritten as

$$S(y + 0) = \max\{T(y - c + 0), S(y) + P(\{y\})\}, \quad (11)$$

$$T(y + 0) = \max\{S(y - c + 0), T(y) + Q(\{y\})\}. \quad (12)$$

We can consider the problem of maximal one to one matching of the points X_i to the points Y_j within the caliper c as the problem of achieving the supremum in (7), where $U(\{(x, y)\}) = k \geq 0$ for $|x - y| \leq c$ means that exactly k points $X_i = x$ are matched to k points $Y_j = y$. Relations (8)–(9) mean that $S(y) - F(y)$ counts the “spare” part of $Q((-\infty, y - c))$, which can not be matched. Analogously $T(y) - G(y)$ counts the part of $P((-\infty, y - c))$ that is not matched. On the other hand, relations (11)–(12) mean that the matching, which corresponds to the measure Z , is done according to the algorithms above. Since we can not match more than $\mu(P, Q)$ pairs, the former algorithm yields the maximal number of matched pairs.

For the case of 1-to- n matching we can put

$$\tilde{P}(\{y\}) = n \cdot \#\{j : X_j = y\}, \quad \tilde{Q}(\{y\}) = \#\{j : Y_j = y\}.$$

and repeat the above argument.

Remark. The first proof of the optimality of the above algorithms can be used for a new proof of the main theorem in Ruzankin (2001).

References

1. Pimentel, S. D., Kelz, R. R., Silber, J. H., and Rosenbaum, P. R. (2015), “Large, Sparse Optimal Matching With Refined Covariate Balance in an Observational Study of the Health Outcomes Produced by New Surgeons,” *Journal of the American Statistical Association*, 110, No. 510, 517–527.
2. Rosenbaum, P. R. (2012), “Optimal Matching of an Optimally Chosen Subset in Observational Studies,” *Journal of Computational and Graphical Statistics*, 21, No. 1, 57–71.
3. Ruzankin, P. S. (2001), “Construction of the optimal joint distribution of two random variables,” *Theory Probab. Appl.*, 46, No.2, 316–334.